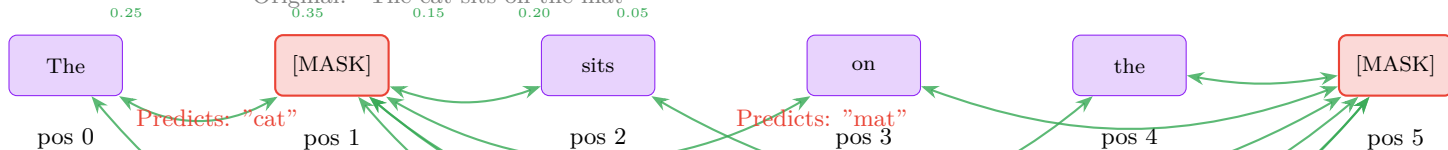


MLM Training Example:

Original: "The cat sits on the mat"
0.35 0.15 0.20 0.05



Bidirectional Attention Matrix:

	The	[MASK]	sits	on	the	[MASK]
The						
[MASK]						

Key Differences from Autoregressive Models

Bidirectional Context: Uses both left & right context
No Causal Masking: Future tokens inform predictions
Parallel Prediction: Multiple masks predicted jointly
MLM Objective: Predict original tokens using full bidirectional context